

Tugas Kelompok

Nama Kelompok

1. Moh Naufal Faqih (10222044)
2. Firman Firdaus (10222033)
3. Ryan Azis Saputra (10222041)
4. Andika Fitra Ramadan (10222043)

Mata Kuliah

Dosen

Program Studi

Pengolahan dan Analisis Sains Data

Dr. Susmini Indriani Lestariningati, S.T, M.T.

Sistem Komputer

## 1. Deteksi Missing Value

### a. Program

```
1 missingValues = []
2
3 for i, row in enumerate(data):
4     for j, value in enumerate(row):
5         if value.strip() == '' or value.lower() == "nan" or value is None:
6             missingValues.append((i + 1, header[j]))
7 if missingValues:
8     print(f"Missing value Ditemukan pada:")
9     for row, col in missingValues:
10        print(f" - Baris {row}, Kolom {col}")
11    print(f"Total Data yang Hilang: {len(missingValues)}")
12 else:
13    print("Tidak Ditemukan Missing Value.")
```

### b. Output

```
Missing value Ditemukan pada:
 - Baris 1, Kolom HARGA
 - Baris 1, Kolom LT
 - Baris 9, Kolom HARGA
 - Baris 19, Kolom HARGA
Total Data yang Hilang: 4
```

## 2. Deteksi Duplikasi Data

### a. Program

```
1  duplikasi_identik = {}
2
3  for i, row in enumerate(data):
4
5      key = tuple(row)
6      if key in duplikasi_identik:
7          duplikasi_identik[key].append(i + 1)
8      else:
9          duplikasi_identik[key] = [i + 1]
10
11
12 duplikasi_identik = {key: value for key, value in duplikasi_identik.items() if len(value) > 1}
13 totalDuplikasiIdentik = sum(len(rows) for rows in duplikasi_identik.values())
14
15
16 if duplikasi_identik:
17     print("Data identik yang terduplikasi ditemukan:")
18     for key, rows in duplikasi_identik.items():
19         print(f" - Data: {key} ditemukan di baris: {rows}")
20     print(f"Total data identik yang terduplikasi: {totalDuplikasiIdentik}")
21 else:
22     print("Tidak ada data identik yang terduplikasi.")
```

### b. Output

```
Data identik yang terduplikasi ditemukan:
- Data: ('7,600,000,000', '278', '350', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [7, 109]
- Data: ('7,000,000,000', '384', '400', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [14, 285]
- Data: ('5,800,000,000', '144', '285', '4', '3', 'ADA', 'JAKSEL') ditemukan di baris: [19, 318, 709]
- Data: ('40,000,000,000', '1500', '700', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [23, 125]
- Data: ('5,680,000,000', '151', '245', '3', '3', 'ADA', 'JAKSEL') ditemukan di baris: [27, 544]
- Data: ('10,500,000,000', '220', '350', '4', '3', 'ADA', 'JAKSEL') ditemukan di baris: [28, 435]
- Data: ('7,050,000,000', '128', '315', '3', '3', 'ADA', 'JAKSEL') ditemukan di baris: [29, 823]
- Data: ('68,000,000,000', '470', '1000', '5', '5', 'ADA', 'JAKSEL') ditemukan di baris: [32, 837]
- Data: ('85,000,000,000', '770', '500', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [34, 844]
- Data: ('20,900,000,000', '754', '500', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [36, 319]
- Data: ('14,000,000,000', '410', '800', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [38, 299]
- Data: ('14,000,000,000', '460', '400', '4', '3', 'ADA', 'JAKSEL') ditemukan di baris: [71, 607]
- Data: ('8,000,000,000', '323', '317', '5', '4', 'ADA', 'JAKSEL') ditemukan di baris: [75, 107]
- Data: ('10,000,000,000', '395', '450', '5', '5', 'ADA', 'JAKSEL') ditemukan di baris: [79, 592]
- Data: ('5,500,000,000', '234', '115', '2', '2', 'ADA', 'JAKSEL') ditemukan di baris: [93, 785]
- Data: ('16,000,000,000', '407', '575', '5', '4', 'ADA', 'JAKSEL') ditemukan di baris: [122, 717]
- Data: ('55,000,000,000', '1500', '1200', '6', '6', 'ADA', 'JAKSEL') ditemukan di baris: [128, 239]
- Data: ('21,000,000,000', '1042', '550', '7', '6', 'ADA', 'JAKSEL') ditemukan di baris: [144, 722]
- Data: ('22,000,000,000', '455', '570', '4', '5', 'ADA', 'JAKSEL') ditemukan di baris: [149, 228, 740]
- Data: ('15,000,000,000', '714', '750', '5', '5', 'ADA', 'JAKSEL') ditemukan di baris: [155, 167]
- Data: ('8,000,000,000', '380', '450', '5', '5', 'ADA', 'JAKSEL') ditemukan di baris: [161, 191]
- Data: ('18,900,000,000', '470', '660', '4', '4', 'TIDAK ADA', 'JAKSEL') ditemukan di baris: [162, 192]
- Data: ('21,900,000,000', '377', '260', '6', '5', 'ADA', 'JAKSEL') ditemukan di baris: [164, 193]
- Data: ('5,600,000,000', '269', '300', '4', '4', 'ADA', 'JAKSEL') ditemukan di baris: [165, 194]
...
- Data: ('700,000,000', '60', '52', '3', '2', 'ADA', 'JAKSEL') ditemukan di baris: [797, 814]
- Data: ('14,900,000,000', '291', '400', '4', '4', 'TIDAK ADA', 'JAKSEL') ditemukan di baris: [802, 857]
- Data: ('1,200,000,000', '60', '70', '3', '2', 'ADA', 'JAKSEL') ditemukan di baris: [806, 818, 890]
Total data identik yang terduplikasi: 130
```

### c. Catatan Tambahan

- Deteksi pada duplikasi hanya dilakukan pada data yang benar-benar identik, dengan mempertimbangkan nilai pada kolom lainnya. Seperti luas tanah (LT), luas bangunan (LB), jumlah kamar tidur (JKT), jumlah kamar mandi (JKM), dan garasi (GRS).
- Jika deteksi duplikasi hanya dilakukan pada nilai kolom harga saja itu akan terdeteksi 735 data yang harganya mempunyai nilai yang sama.

## 3. Deteksi Outlier Menggunakan Metode IQR

### a. Program

```
1 # Deteksi Outlier
2 def hitungIQR(data):
3     data.sort()
4     n = len(data)
5     q1 = data[n // 4] if n % 4 == 0 else (data[n // 4] + data[n // 4 - 1]) / 2
6     q3 = data[3 * n // 4] if (3*n) % 4 == 0 else (data[3 * n // 4] + data[3*n // 4 - 1]) / 2
7     iqr = q3 - q1
8     return q1, q3, iqr
9
10 dataHarga = [float(row[0].replace(',', '')) for row in data if row[0].strip() != '' and row[0].lower() != 'nan']
11 q1, q3, iqr = hitungIQR(dataHarga)
12
13 lowerBound = q1 - 1.5 * iqr
14 upperBound = q3 + 1.5 * iqr
15
16 outliers = [value for value in dataHarga if value < lowerBound or value > upperBound]
17 totalOutliers = len(outliers)
18 print(f"Total Outlier: {totalOutliers}")
19
20 print(f"Q1: {q1}, Q3: {q3}, IQR: {iqr}")
21 print(f"batas bawah: {lowerBound}, batas atas: {upperBound}")
22 if outliers:
23     print(f"Outlier ditemukan: {outliers}")
24 else:
25     print("Tidak ada outlier ditemukan")
```

### b. Output

Total Outlier: 65

Q1: 6750000000.0, Q3: 20000000000.0, IQR: 13250000000.0

batas bawah: -13125000000.0, batas atas: 39875000000.0

Outlier ditemukan: [40000000000.0, 40000000000.0, 40000000000.0, 40000000000.0, 40000000000.0, 42000000000.0, 42000000000.0, 42000000000.0, 43000000000.0, 43000000000.0, 45000000000.0, 45000000000.0, 45000000000.0, 47000000000.0, 49000000000.0, 50000000000.0, 50000000000.0, 55000000000.0, 55000000000.0, 55000000000.0, 55000000000.0, 55000000000.0, 55000000000.0, 56000000000.0, 57000000000.0, 60000000000.0, 60000000000.0, 60000000000.0, 60000000000.0, 65000000000.0, 65000000000.0, 68000000000.0, 68000000000.0, 68000000000.0, 72000000000.0, 75000000000.0,

75000000000.0, 75000000000.0, 77500000000.0, 80000000000.0,  
 85000000000.0, 85000000000.0, 85000000000.0, 85000000000.0,  
 85500000000.0, 90000000000.0, 95000000000.0, 100000000000.0,  
 110000000000.0, 120000000000.0, 135000000000.0, 165000000000.0,  
 168000000000.0, 169000000000.0, 175230000000.0, 180000000000.0,  
 185000000000.0, 250000000000.0]

**c. Catatan Tambahan:**

- Total outlier sebanyak 65 data terdeteksi sebagai outlier.
- Ini menunjukkan bahwa ada banyak data yang berada di luar batas distribusi normal harga rumah.
- Q1 (Kuartil Pertama) 6.75 miliar (6.750.000.000.0) menunjukkan 25% dari data memiliki harga rumah di bawah nilai ini.
- Q3 (Kuartil Ketiga) 20 miliar (20.000.000.000.0) menunjukkan 75% dari data memiliki harga rumah di bawah nilai ini.
- IQR (Interquartile Range) 13.25 miliar (13.250.000.000.0) menunjukkan rentang variasi harga rumah antara Q1 dan Q3 di tengah distribusi.
- Harga rumah yang sangat tinggi (di atas 39.875 miliar) kemungkinan besar mencerminkan properti mewah atau lokasi premium dengan mempertimbangkan faktor faktor lainnya, seperti luas tanah, luas bangunan, jumlah kamar tidur, jumlah kamar mandi, dan garasi.

## 4. Proses Menangani Missing Value

**a. Program**

```
1 cleaned_data = [row for row in data if all(value.strip() != "" and value.lower() != "nan" and value is not None for value in row)]
2
3 # Output hasil
4 print(f"Jumlah data sebelum pembersihan: {len(data)}")
5 print(f"Jumlah data setelah pembersihan: {len(cleaned_data)}")
6 print(f"Jumlah baris yang dihapus: {len(data) - len(cleaned_data)}")
```

**b. Output**

```
Jumlah data sebelum pembersihan: 1001
Jumlah data setelah pembersihan: 998
Jumlah baris yang dihapus: 3
```

**c. Catatan Tambahan**

- Karena data yang missing/hilang hanya sedikit yaitu 3, melakukan penghapusan pada data adalah keputusan yang tepat dibanding melakukan mengisi (imputasi) yang bisa menambah bias atau error.

- Karena data hanya 0.3% yang dihapus dari seluruh data maka tidak akan memengaruhi rata-rata, median maupun sebaran data.

## 5. Proses Menangani Duplikasi Data

### a. Program

```

1 dataUnik = []
2 seen = set()
3
4 for row in cleaned_data:
5     rowTuple = tuple(row)
6     if rowTuple not in seen:
7         seen.add(rowTuple)
8         dataUnik.append(row)
9 print(f"Jumlah data sebelum penghapusan duplikasi: {len(cleaned_data)}")
10 print(f"Jumlah data setelah penghapusan duplikasi: {len(dataUnik)}")
11 print(f"Jumlah baris yang dihapus: {len(cleaned_data) - len(dataUnik)}")
12

```

### b. Output

```

Jumlah data sebelum penghapusan duplikasi: 998
Jumlah data setelah penghapusan duplikasi: 929
Jumlah baris yang dihapus: 69

```

## 6. Proses Menangani Outlier

### a. Program

```

1 dataHarga = [float(row[0].replace(',', '')) for row in dataUnik if row[0].strip() != '' and row[0].lower() != 'nan']
2
3 # Hitung IQR
4 q1, q3, iqr = hitungIQR(dataHarga)
5 lowerBound = q1 - 1.5 * iqr
6 upperBound = q3 + 1.5 * iqr
7
8 segmentasiData = []
9 for row in dataUnik:
10     harga = float(row[0].replace(',', ''))
11     if harga > lowerBound and harga < iqr:
12         row.append("Harga Rumah di Bawah Harga Reguler")
13     elif harga > iqr and harga < upperBound:
14         row.append("Rumah Reguler")
15     elif harga > upperBound:
16         row.append("Rumah Mewah")
17     segmentasiData.append(row)
18 print(f"Jumlah rumah di bawah Harga Reguler: {sum(1 for row in segmentasiData if row[-1] == 'Harga Rumah di Bawah Harga Reguler')}")
19 print(f"Jumlah rumah reguler: {sum(1 for row in segmentasiData if row[-1] == 'Rumah Reguler')}")
20 print(f"Jumlah rumah mewah: {sum(1 for row in segmentasiData if row[-1] == 'Rumah Mewah')}")
21

```

## b. Output

```
Jumlah rumah di bawah Harga Reguler: 467
Jumlah rumah reguler: 402
Jumlah rumah mewah: 60
```

## c. Catatan Tambahan

- Outlier pada data ditandai dengan membaginya ke beberapa segmen berdasarkan harganya (misalnya: di bawah harga reguler dan mewah).
- Ada beberapa rumah dengan harga sangat tinggi yang kemungkinan memiliki fasilitas eksklusif dan lengkap.
- Outlier tidak dianggap sebagai “kesalahan” atau “anomali”, tapi dijadikan kategori khusus yang memberi makna baru dalam analisis.
- Dengan mengelompokkan harga ekstrem ke segmen tersendiri, dapat menghindari distorsi pada analisis rata-rata, atau median, karena outlier tak lagi mempengaruhi data utama

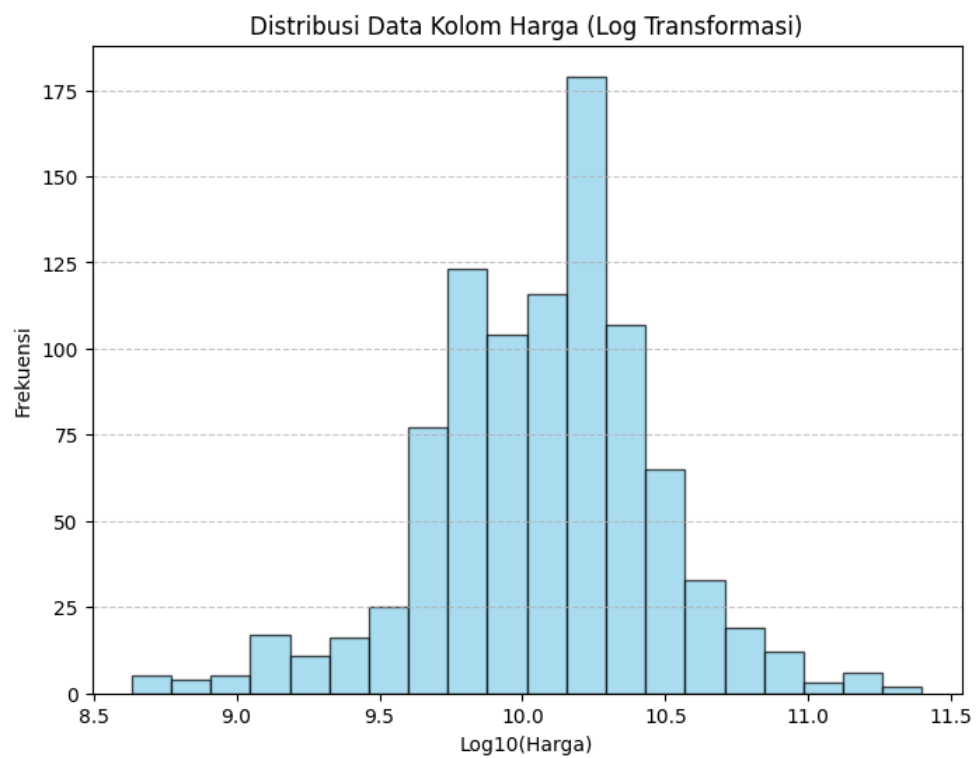
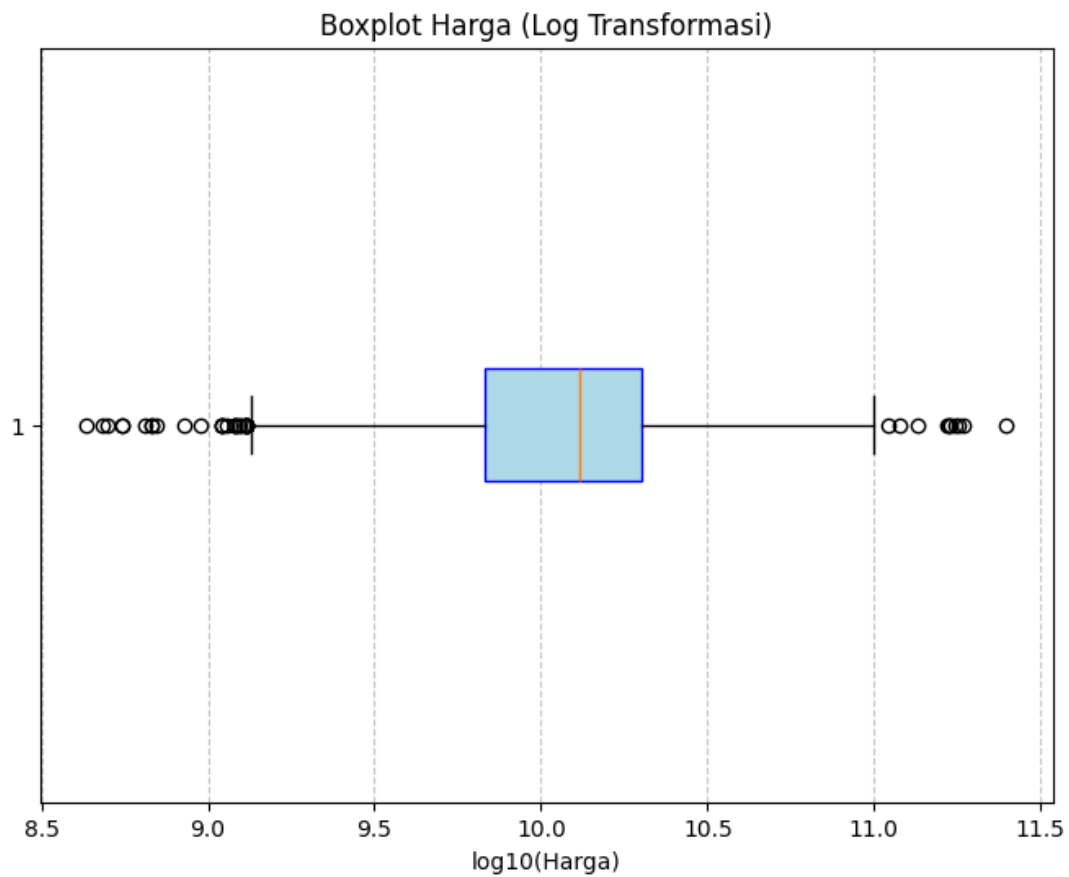
## 7. Visualisasi Data (Boxplot dan Histogram)

### a. Program

```
1 # visualisasi data( boxplot)
2 log_data_harga = np.log10(dataHarga)
3 plt.figure(figsize=(8, 6))
4 plt.boxplot(log_data_harga, vert=False, patch_artist=True, boxprops=dict(facecolor='lightblue', color='blue'))
5 plt.title("Boxplot Harga (Log Transformasi)")
6 plt.xlabel("log10(Harga)")
7 plt.grid(axis='x', linestyle='--', alpha=0.7)
8
9 # Tampilkan plot
10 plt.show()
```

```
1 log_data_harga = np.log10(dataHarga)
2
3 # Membuat histogram
4 plt.figure(figsize=(8, 6))
5 plt.hist(log_data_harga, bins=20, color='skyblue', edgecolor='black', alpha=0.7)
6 plt.title("Distribusi Data Kolom Harga (Log Transformasi)")
7 plt.xlabel("Log10(Harga)")
8 plt.ylabel("Frekuensi")
9 plt.grid(axis='y', linestyle='--', alpha=0.7)
10
11 # Tampilkan plot
12 plt.show()
```

## b. Output



### **c. Catatan Tambahan**

- Nilai pada harga rumah awalnya condong ke kiri, artinya Sebagian besar harga ada di level rendah-menengah tapi ada segelintir yang sangat mahal.
- Maka dari itu dilakukanlah transformasi logaritmik (misal:  $\log_{10}(\text{harga})$ ) menormalkan distribusi, membuatnya lebih mendekati distribusi normal sehingga lebih mudah dianalisis secara statistik.
- Tanpa log transform, data mahal bisa menutupi variasi data di harga rendah.
- Dengan log, semua nilai bisa terlihat lebih proporsional, sehingga struktur sebaran data menjadi lebih jelas (seperti kelompok harga tertentu atau outlier)
- Log transformasi menekan nilai-nilai ekstrem (contohnya rumah 200 miliar) sehingga outlier tidak terlalu mendistorsi skala grafik.